

Grid Developments & Operations

for ATLAS Experiment on Behalf of IT Cloud

Manoj Kumar Jha
INFN- CNAF, Bologna

23rd Jan., 2012

- Development of grid tools
 - Ganga: User friendly job submission and management tool
 - Functional test with GangaRobot
 - ATLAS task book keeping
- Grid operations
 - Tier0 data registered and exported
 - Overview of problem
 - Data distribution
 - Storage
 - Software performance
- New Ideas !
- Other activities

Data Analysis with Ganga

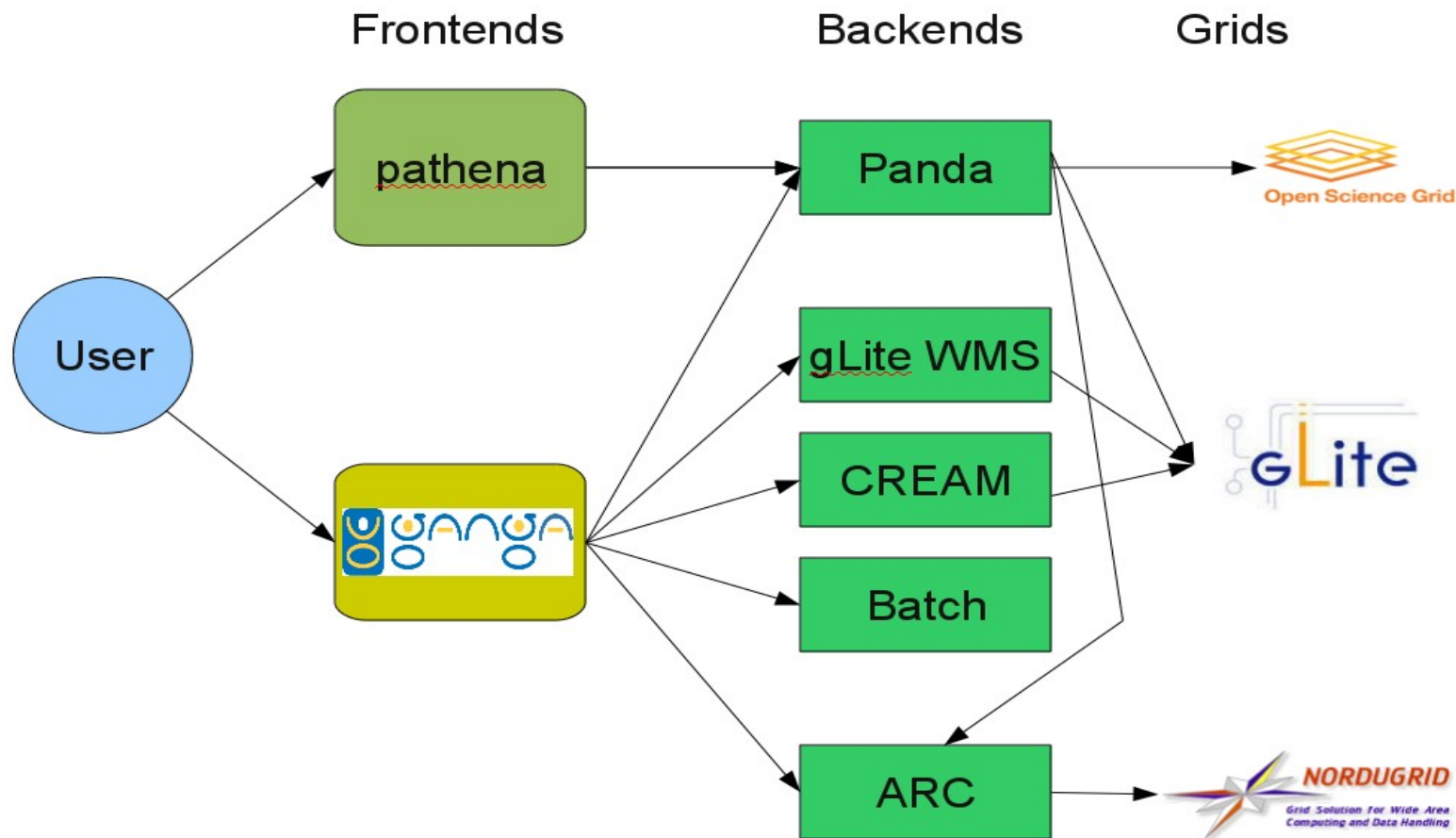
Accepted for publication in J. Phys. Conf. Series



Challenges in a LHC Data Analysis

- Data volumes
 - LHC experiments produce and store several PetaBytes /year
 - ATLAS recorded ~ 5.2 fb⁻¹ of data till now
- CPUs
 - Event complexity and number of users demands: at least 100000 CPUs based on computing model
- Software
 - The experiments have complex software environment and framework
- Connectivity
 - Data should be available at 24/7 at a high bandwidth
- Distributed analysis tools must should be
 - Easy to configure and fast to work with
 - Reliable and jobs should have 100% success rate at 1st attempt

Atlas Distributed Analysis Layers



Data is centrally being distributed by DQ2 – Jobs go to data

Introduction to Ganga

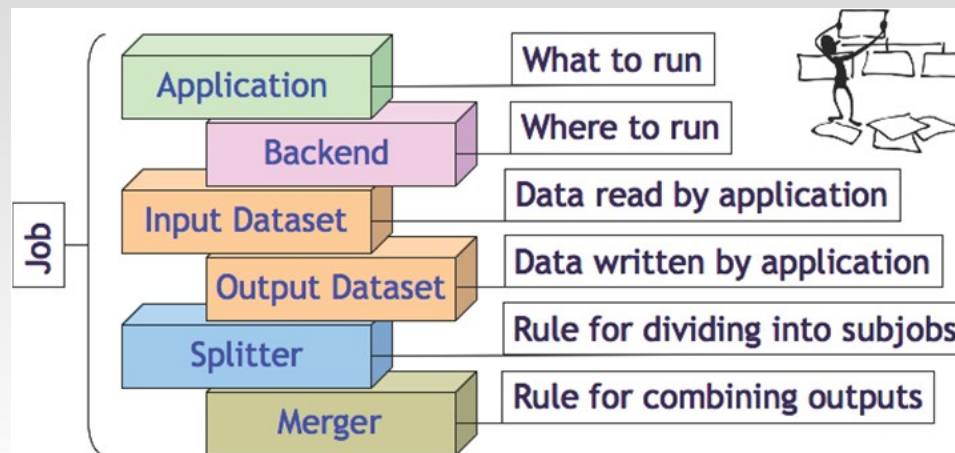


- **Ganga is a user-friendly job management tool.**
 - Jobs can run locally or on a number of batch systems and grids.
 - Easily monitor the status of jobs running everywhere.
 - To change where the jobs run, change one option and resubmit.
- **Ganga is the main distributed analysis tool for LHCb and ATLAS.**
 - Experiment-specific plugins are included.
- **Ganga is an open source community-driven project:**
 - Core development is joint between LHCb and ATLAS
 - Modular architecture makes it extensible by anyone
 - Mature and stable, with an organized development process

Submitting a Job with Ganga



What is a Ganga Job?



Run the default job locally:

```
Job().submit()
```

Default job on the EGEE grid:

```
Job(backend=LCG()).submit()
```

Listing of the existing jobs:

```
jobs
```

Get help (e.g. on a job):

```
help(jobs)
```

Display the nth job:

```
jobs(n)
```

Copy and resubmit the nth job:

```
jobs(n).copy().submit()
```

Copy and submit to another grid:

```
j=jobs(n).copy()
```

```
j.backend=DIRAC()
```

```
j.submit()
```

Kill and remove the nth job:

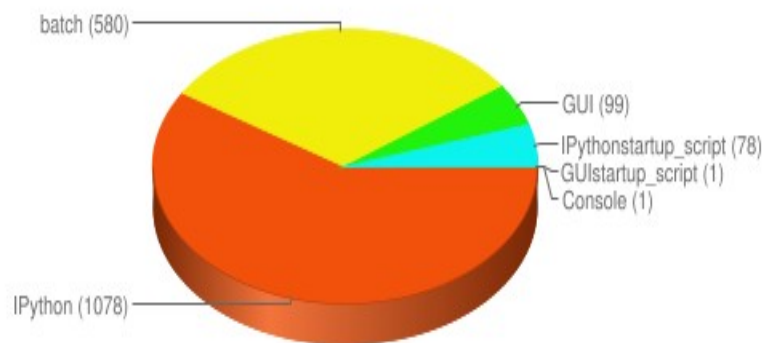
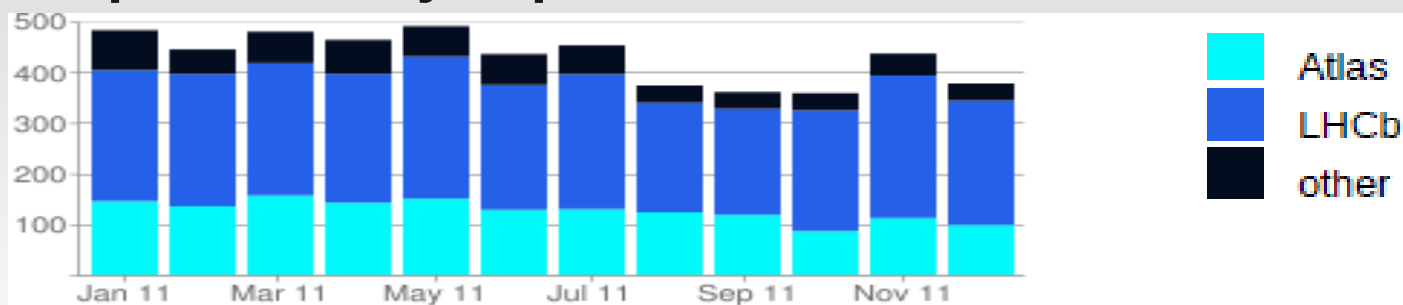
```
job(n).kill()
```

```
job(n).remove()
```

Number of Ganga Users



Unique users by experiment in 2011



- Total number sessions: 364112 Number of unique users: 1107
- Number of sites: 127
- Python scripting is more popular than using Ganga in batch mode.
- GUI is not used often ..., good for tutorials and learning.

Conclusions



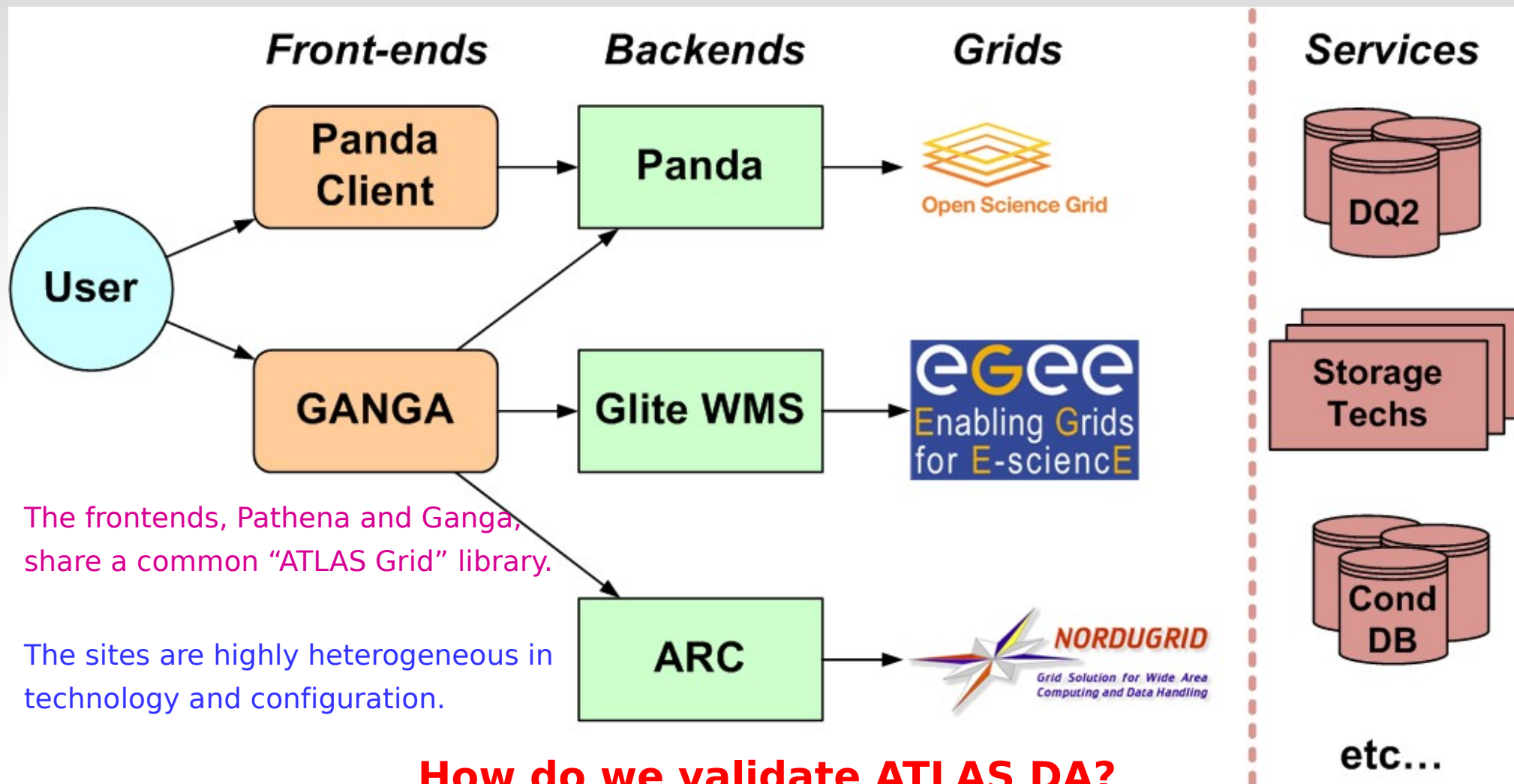
- Ganga is a user-friendly job management tool for Grid, Batch and Local systems
 - “configure once, run anywhere”
- A stable development model:
 - Well organized release procedure with extensive testing
 - Plugin architecture allows new functionality to come from non-core developers
 - Not just a UI – provides a Grid API on which many applications are built
 - Strong development support from LHCb and ATLAS, and 25% usage in other VOs

For more information visit <http://cern.ch/ganga>

Functional Testing with GangaRobot

Accepted for publication in J. Phys. Conf. Series

DA in ATLAS: What are the resources?



How do we validate ATLAS DA?

Use case functionalities?? Behaviour under load??

Functional Testing with GangaRobot



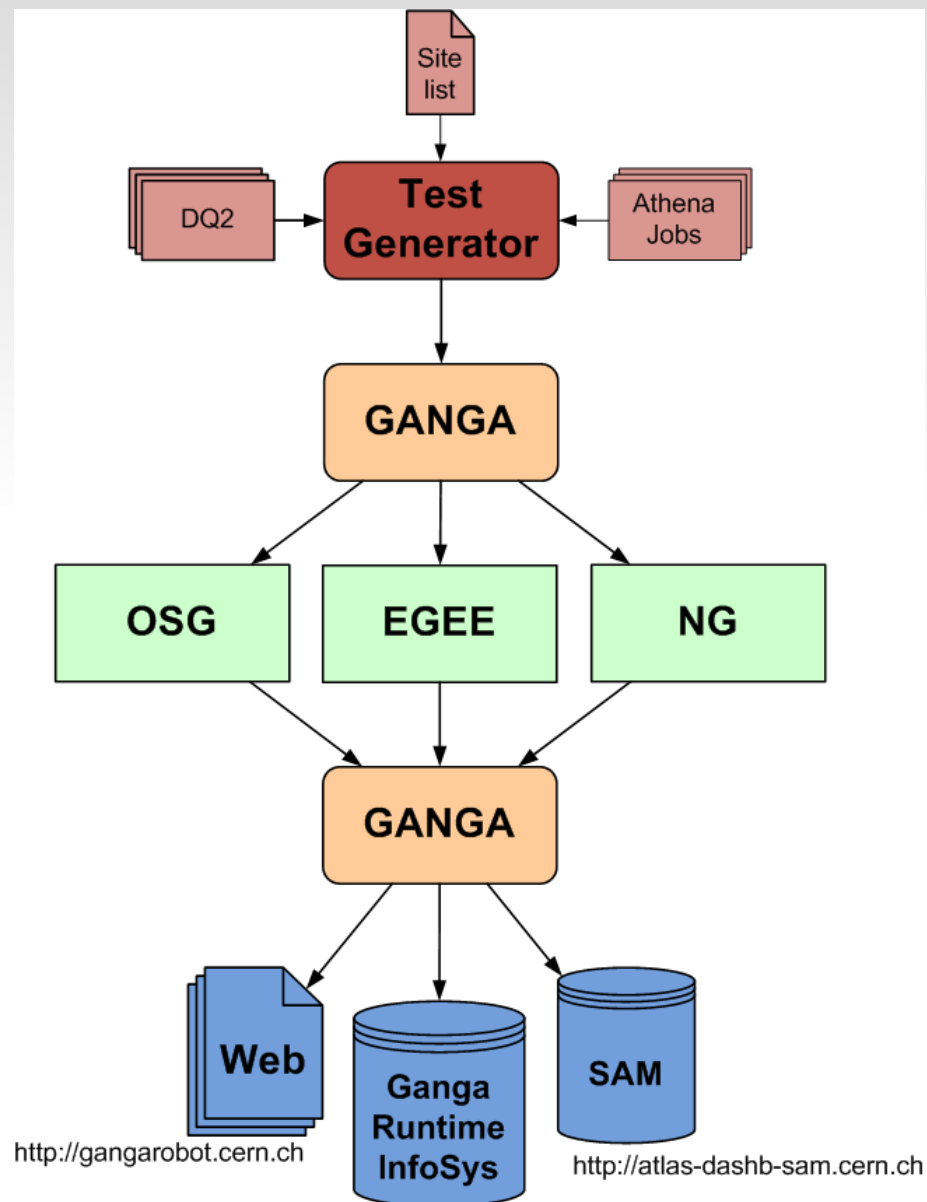
- Definitions:
 - **Ganga** is a distributed analysis user interface with a scriptable python API.
 - **GangaRobot** is both
 - a) a component of Ganga which allows for rapid definition and execution of test jobs, with hooks for pre- and post-processing
 - b) an ATLAS service which uses (a) to run DA functional tests
- So what does GangaRobot test and how does it work?



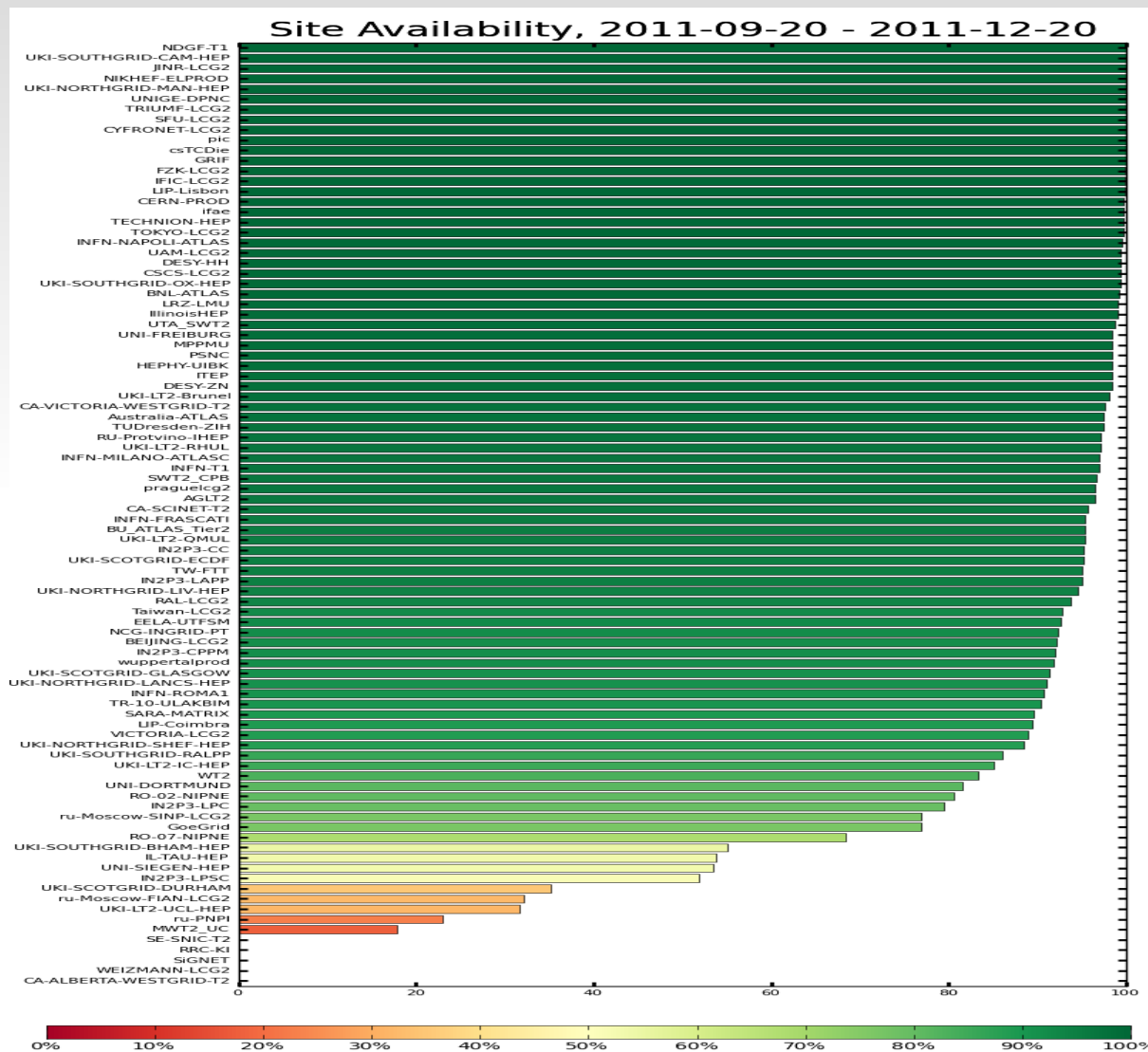
Functional Testing with GangaRobot



1. Tests are defined by the GR operator:
 - Athena version, analysis code, input datasets, which sites to test
 - Short jobs, mainly to test the software and data access
1. Ganga submits the jobs
 - To OSG/Panda, EGEE/LCG, NG/ARC
1. Ganga periodically monitors the jobs until they have completed or failed
 - Results are recorded locally
1. GangaRobot then publishes the results to three systems:
 - Ganga Runtime Info System, to avoid failing sites
 - SAM, so that sites can see the failures
 - GangaRobot website, monitored by ATLAS DA shifters
 - GGUS and RT tickets sent for failures



Overall Statistics with GangaRobot



Plots from SAM dashboard
<http://dashb-atlas-sam.cern.ch/>
 of daily and percentage
 availability of ATLAS sites over
 the past 3 months.

The good: Many sites with
 >90% efficiency

The bad: Some of the sites have
 uptime < 80%

The expected: Many
 transient errors, 1-2 day
 downtimes. A few sites are
 permanently failing.

Conclusions

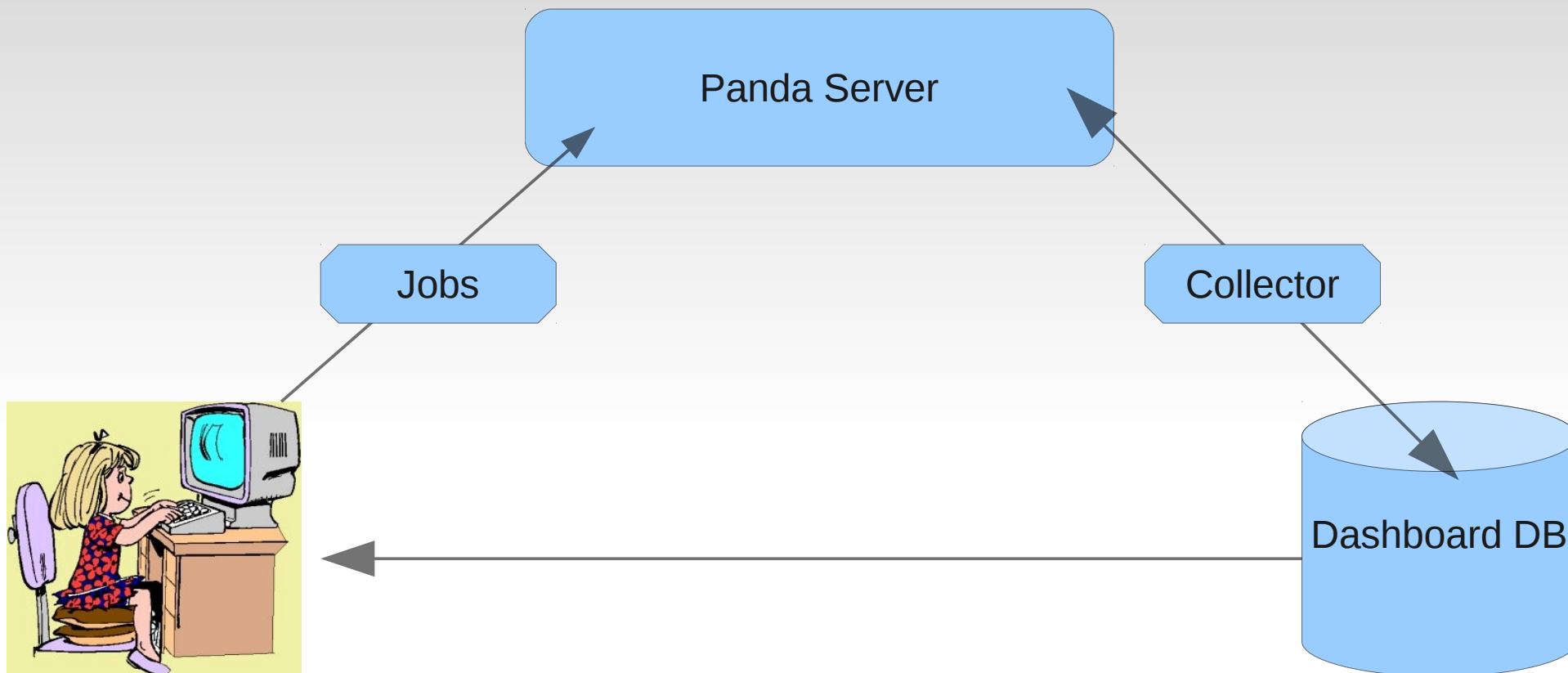


- Validating the grid for user analysis is a top priority for ATLAS Distributed Computing
 - The functionalities available to users are rather complete, now we are testing to see what breaks under full load.
- GangaRobot is an effective tool for functional testing:
 - Daily tests of the common use cases are essential if we want to keep sites working.

ATLAS Task Book Keeping

Under Development

- Analysis job comprises of several subjobs and their associated retried jobs at different sites.
 - All the subjobs belong to same output container dataset, known as task.
- Task API provides
 - Bookkeeping at task level.
 - Information about latest retried jobs
 - Information about number of processed events, files
 - Present a brief summary about task
- Reduce load on PandaDB server by using Dashboard DB.



- ♦ A collector runs at fixed interval of time for getting information from Panda DB and populates it into Dashboard DB. Due to this, there is some latency involved in updating information in dashboard DB with respect to Panda DB (~5 minutes or less) .
- ♦ Executing following url gives information in python object for task 'yourtask' .
 - ♦ <http://dashb-atlas-job.cern.ch/dashboard/request.py/bookkeeping?taskname=yourtask>

```

Terminal
File Edit View Terminal Help
[jha@pcardas02 ~]$ ./AfsHome/public/Atlas/Panda/latest/panda_task.py --outDS user.gabrown.20111017202747.189/

Task name : user.gabrown.20111017202747.189/
Task submission time : 2011-10-17 19:29:24
Task updated time : 2011-10-18 15:56:29
Submission/Modification host : pc194.hep.manchester.ac.uk
Input dataset : data10_7TeV.00167607.physics_JetTauEtmiss.merge.AOD.r1774_p327_p333_tid207070_00
ComputingSite : ANALY_RAL ANALY_CSCS ANALY_QMUL ANALY_LANCS ANALY_INFN-T1
Number of jobs : 195
Status : FINISHED: 193 FAILED: 2
Number of files processed : 1690
Number of events processed : 22256737
JobsetID : 7366 7374
JobDefinitionId : 7371 7368 7370 7369 7367 7375

[jha@pcardas02 ~]$
[jha@pcardas02 ~]$
[jha@pcardas02 ~]$ ./AfsHome/public/Atlas/Panda/latest/panda_task.py --outDS user.gabrown.20111017202747.189/ --status failed --showPandaID

Task name : user.gabrown.20111017202747.189/
Task submission time : 2011-10-17 19:29:24
Task updated time : 2011-10-18 15:56:29
Submission/Modification host : pc194.hep.manchester.ac.uk
Input dataset : data10_7TeV.00167607.physics_JetTauEtmiss.merge.AOD.r1774_p327_p333_tid207070_00
ComputingSite : ANALY_QMUL
Number of jobs : 2
Status : Failed
Number of files processed : 0
Number of events processed : 0
JobsetID : 7374
JobDefinitionId : 7375
Panda IDs : 1336995383-1336995384

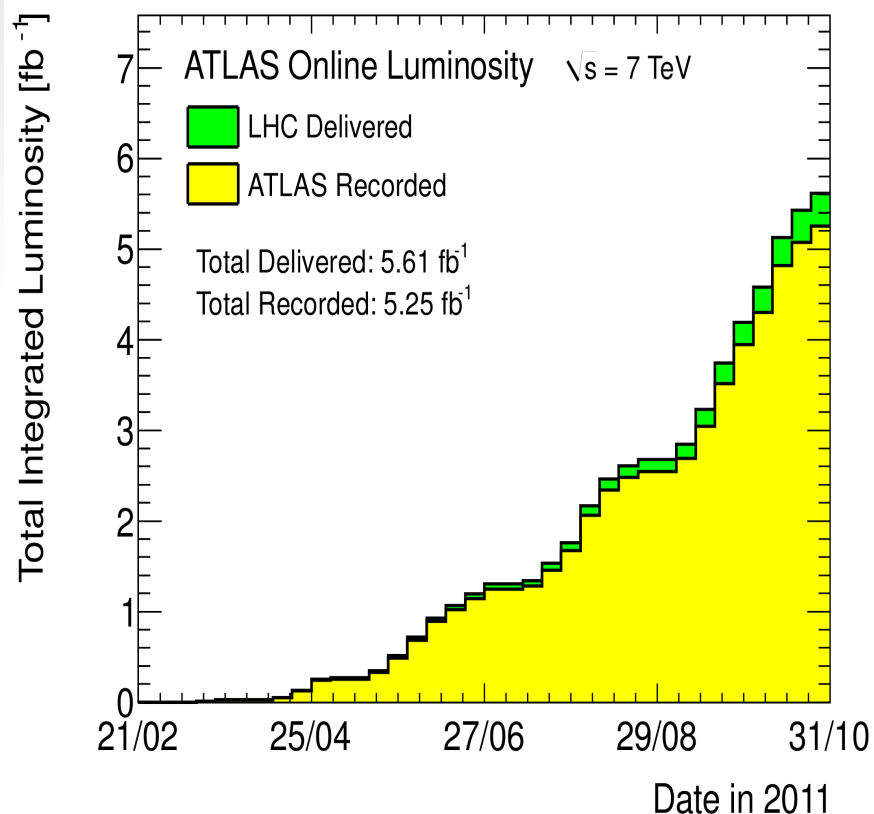
[jha@pcardas02 ~]$

```

- ◆ Task represented by outDS 'user.gabrown.20111017202747.189/ '
- ◆ Total number of jobs: 195
- ◆ Processed at 5 different queues
- ◆ Status : FINISHED: 193 FAILED: 2
- ◆ Second command shows detail information about all the failed jobs.

Grid Operations for ATLAS experiment on behalf of IT Cloud

Introduction: Atlas in Data Taking



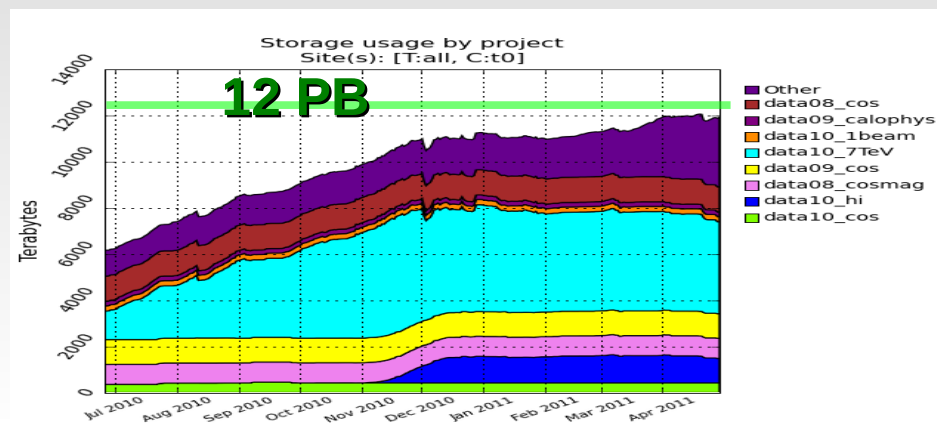
- LHC has been delivering stable beams since 30/03/10.
- ATLAS has been taking data with good efficiency.

Tier-0 Data Registered and Exported

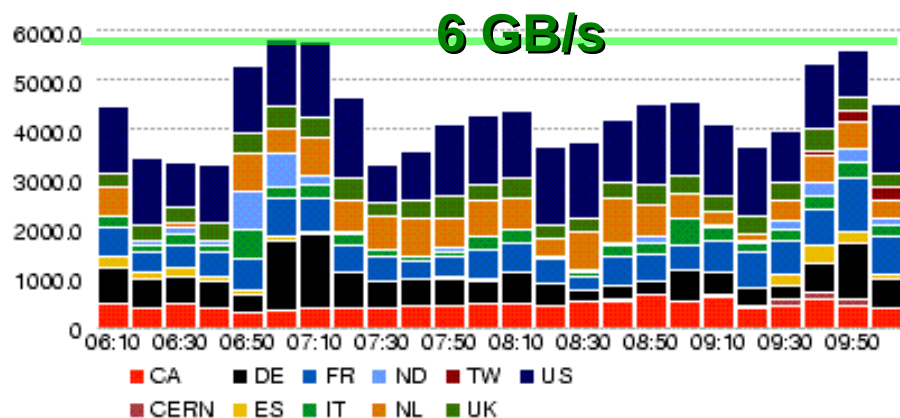


- Data volume registered at Tier-0 since data taking reaching 12 PB
- Data export rate from Tier-0 is more than 5 GB/s
- Some times we need to throttle the export rate in accordance with the available bandwidth at Tier-0

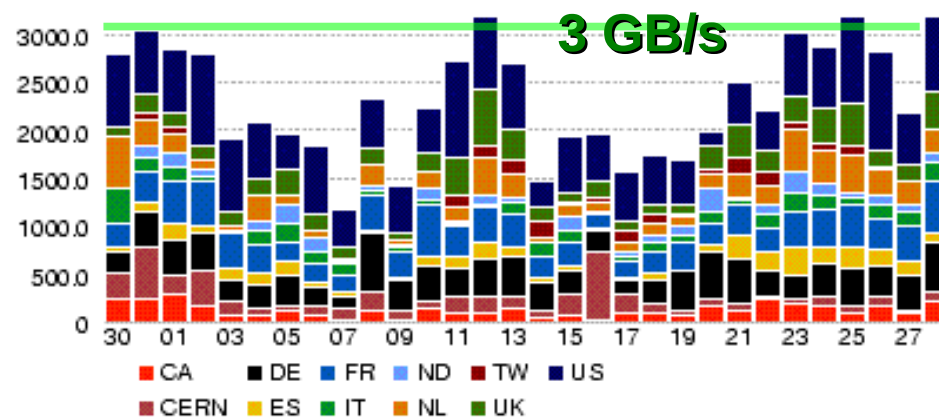
Cumulative data volume registered at Tier-0



Tier-0 export rate: hourly average



Tier-0 export rate: daily average



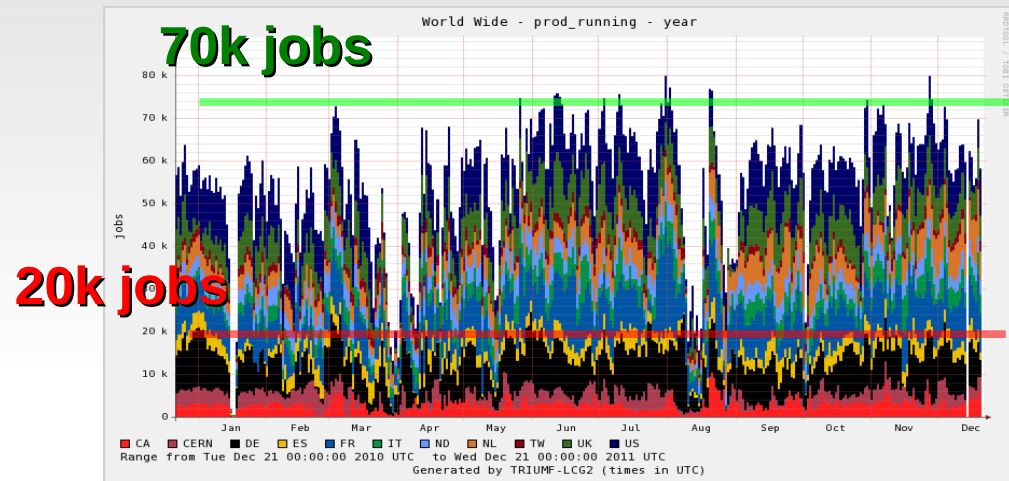
Data Processing Activities



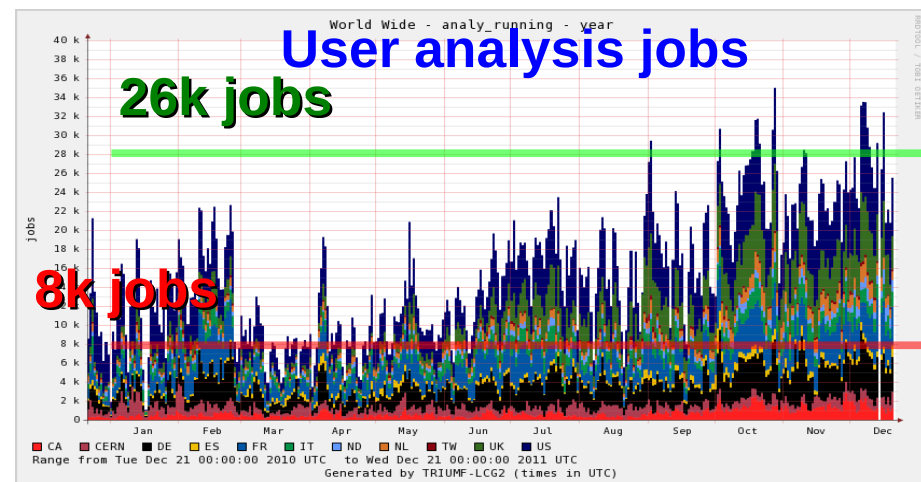
- ATLAS has been able to sustain continued high rate of official production jobs
- Large increase in user analysis jobs
 - The system continues to scale up well

Despite the overall good performance of ATLAS distributed computing, there are bottlenecks available in the system, which we are mentioning in the next slides.

Official production jobs



User analysis jobs



- Distribution Policy
 - Distribution of data using dataset popularity (and unpopularity)
 - Unbalanced data distribution between Tiers
 - Keeping the above factors in mind, it motivates Panda Dynamic Data Placement (PD2PM)
- File corruption
 - File is corrupted using transfer
 - File is corrupted/lost on site
- Communication with user
 - Is the current number of replicas sufficient ?
 - Reconstruction AOD & merged AOD datasets
 - Delay with AOD merging tasks submission lead to many requests for the reconstruction AOD datasets transfer
 - Dataset container content

Overview of Problem: Storage



- Storage instability
 - Storage availability has increased in last years but users
 - Expect job reliability of 100%
 - Still more important than processing speed
- Files with bad checksums
 - Discovered by users/reprocessing jobs (few files per month)
- Lost files
 - It is necessary to have 2 copies of very important data
- Deletion service
 - Sometime files on storage element are not deleted: SE or deletion issue

Overview of Problem: Software Performance

- Growth of static memory squeezes breathing space of Event Data Model
- With increase trigger and pileup rate, CPU/memory usage is going to increase in coming days
 - How to reduce it ?
 - Since a large part of memory used is static, share memory between reconstruction jobs: Athena Multi Process (AthenaMP)

New Ideas !

Data Distribution

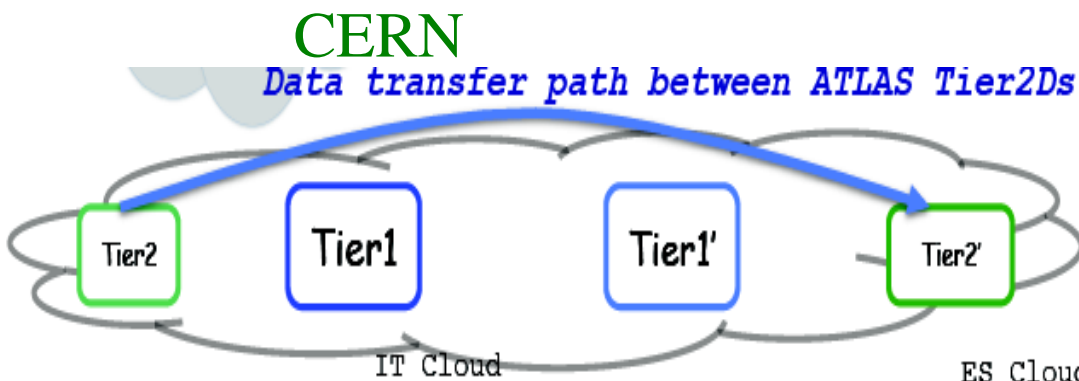


- Local File Catalogs consolidation
 - There are more than 15 LFC ATLAS wide, roughly one per cloud + 6 catalogs in US. If LFC is down, the whole cloud is down. It will be one catalog at CERN and hot backup in another geographical location
- PD2P: Panda Dynamic Data Placement
 - Analysis jobs triggers replication of input data to another site
- T2Ds: Directly connected Tier2
 - Tier2 with the direct connection to CERN

T2Ds commission and sonar test results

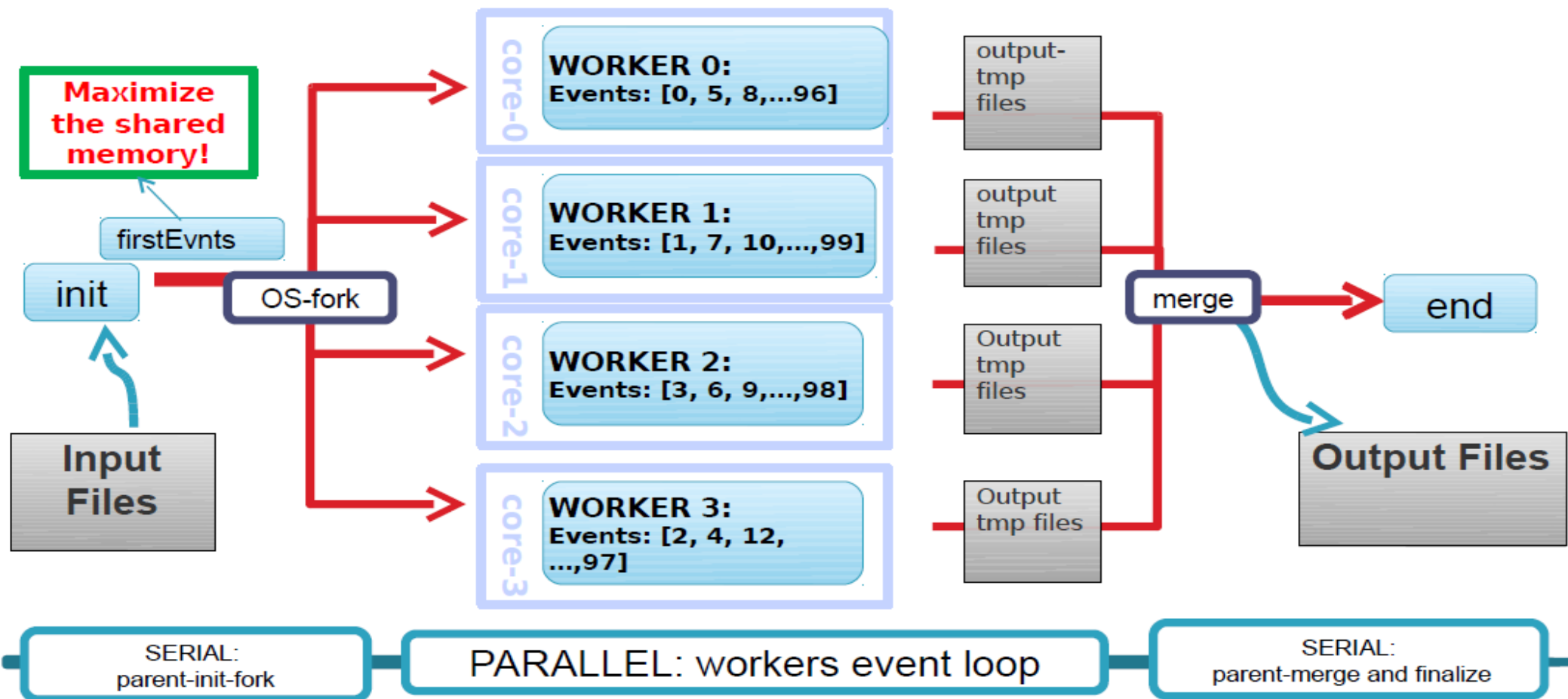
<http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview?view=Sonar>

Site Name	SrcSite	SrcCloud	SrcTier	DstSite	DstCloud	DstTier	AvgBRS(MB/s)	EvS	AvgBRM(MB/s)	EvM	AvgBRL(MB/s)	EvL
AGLT2 to TRIUMF-LCG2	AGLT2	US	T2D	TRIUMF-LCG2	CA	T1	0.86±0.12	10	2.13±0.28	10	4.80±1.29	10
AGLT2 to UAM-LCG2	AGLT2	US	T2D	UAM-LCG2	ES	T2D	0.49±0.18	10	3.10±0.74	10	11.95±4.09	10
AGLT2 to UK-LT2-QMUL	AGLT2	US	T2D	UK-LT2-QMUL	UK	T2D	0.64±0.09	10	3.82±0.92	10	15.47±4.32	5
AGLT2 to UK-LT2-RHUL	AGLT2	US	T2D	UK-LT2-RHUL	UK	T2	0.53±0.07	5	2.03±0.12	10	0.00±0.00	0
AGLT2 to UK-LT2-UCL-HEP	AGLT2	US	T2D	UK-LT2-UCL-HEP	UK	T2	0.44±0.06	5	1.45±0.22	5	0.00±0.00	0
AGLT2 to UK-NORTHGRID-LANCS-HEP	AGLT2	US	T2D	UK-NORTHGRID-LANCS-HEP	UK	T2D	0.66±0.03	10	4.25±1.00	10	21.22±3.87	5
AGLT2 to UK-NORTHGRID-LIV-HEP	AGLT2	US	T2D	UK-NORTHGRID-LIV-HEP	UK	T2	0.49±0.05	5	1.58±0.18	10	0.00±0.00	0
AGLT2 to UK-NORTHGRID-MAN-HEP	AGLT2	US	T2D	UK-NORTHGRID-MAN-HEP	UK	T2D	0.55±0.07	10	3.50±1.37	10	18.42±9.72	10
AGLT2 to UK-NORTHGRID-SHEF-HEP	AGLT2	US	T2D	UK-NORTHGRID-SHEF-HEP	UK	T2	0.55±0.08	5	1.91±0.60	10	0.00±0.00	0
AGLT2 to UK-NORTHGRID-SCOTGRID-ECDE	AGLT2	US	T2D	UK-NORTHGRID-SCOTGRID-ECDE	UK	T2	0.53±0.04	5	4.49±0.15	5	0.00±0.00	0
AGLT2 to UK-NORTHGRID-GLASGOW	AGLT2	US	T2D	UK-NORTHGRID-GLASGOW	UK	T2D	0.54±0.05	10	1.45±0.85	10	9.48±3.44	7
AGLT2 to UK-NORTHGRID-BHAM-HEP	AGLT2	US	T2D	UK-NORTHGRID-BHAM-HEP	UK	T2	0.58±0.04	5	2.80±0.53	10	0.00±0.00	0
AGLT2 to UK-NORTHGRID-CAM-HEP	AGLT2	US	T2D	UK-NORTHGRID-CAM-HEP	UK	T2	0.43±0.11	5	2.36±0.47	10	0.00±0.00	0
AGLT2 to UK-NORTHGRID-OR-HEP	AGLT2	US	T2D	UK-NORTHGRID-OR-HEP	UK	T2	0.63±0.07	10	2.69±1.14	10	0.00±0.00	0
AGLT2 to UK-NORTHGRID-BALPP	AGLT2	US	T2D	UK-NORTHGRID-BALPP	UK	T2	0.30±0.07	5	0.56±0.41	10	0.00±0.00	0



Event Level Parallelism with AthenaMP

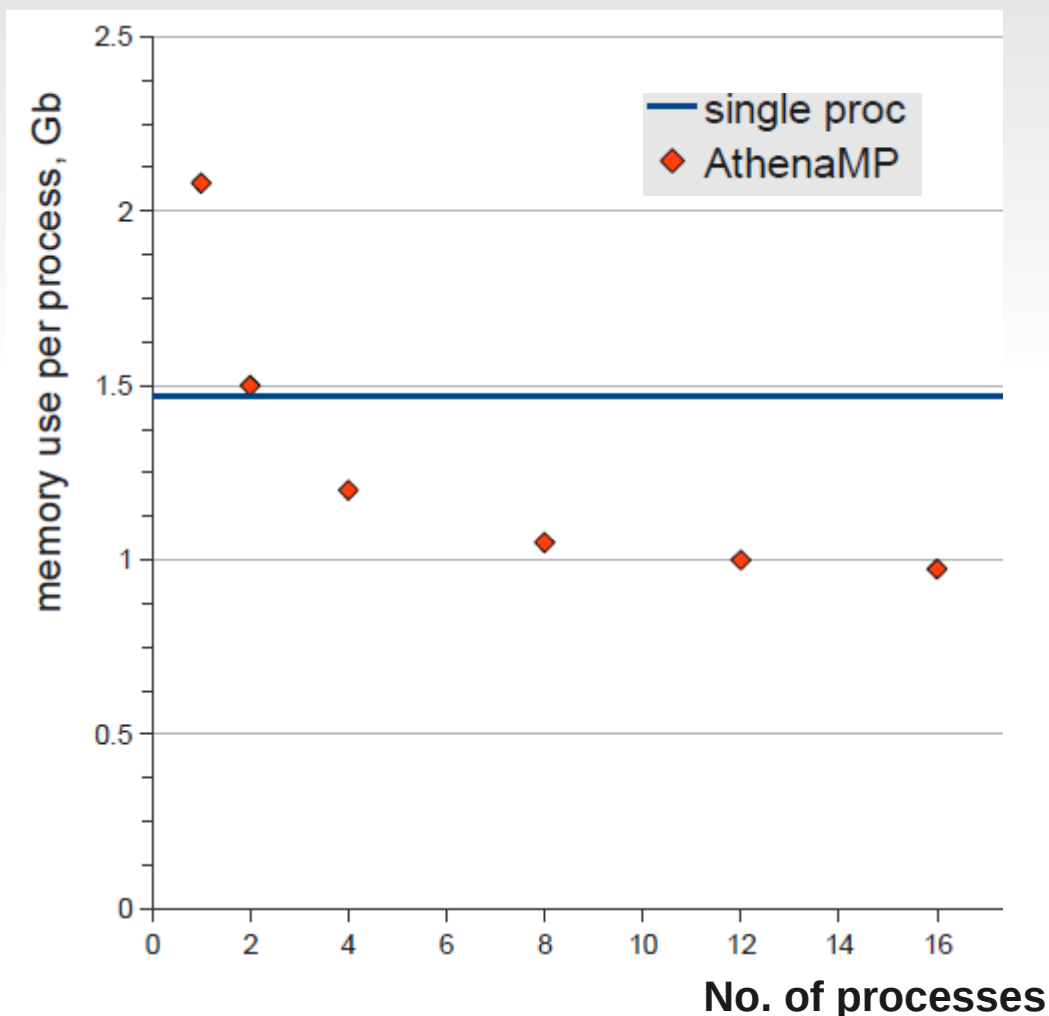
```
> Athena.py --nprocs=4 -c EvtMax=100 Jobo.py
```



Why AthenaMP ?



8 core HT machine



- Main goal is to reduce overall memory footprint
- Use linux fork() to share memory automatically
- **AthenaMP ~0.5 Gb physical memory saved per process**

Looking to the Future



- Beyond dynamic data placement
 - Event level caching
- Cloud computing
 - Investigation of “Amazon S3” or similar web based protocols to access/integrate cloud storage in the medium/longer term.
- Highly scalable 'noSQL' database (it is not replacement of ORACLE, but most probably we will have a hybrid of two technologies).
- Monitoring, diagnostics, error management automation
- CERNVM: Portable analysis environment using actualization technology

Key Issues for ADC in 2012



- Maintain reliable and robust MC production and data reprocessing over grid
- Full support of physics group production
- Reliable access to data ATLAS wide
 - Minimize possibility of single point failure
 - Commissioning Tier2Ds
- Distributed analysis
 - User's support
 - Distributed analysis back end and front end unification
 - Evolution of user's support interface, like web based support forum which complements egroup

Other Activities

- Did PhD on CMS experiment in March 2007.
 - Proposed the geometry of lead absorbers in Preshower of CMS detector. This geometry was accepted by CMS ECAL group.
- Visiting Scientist at LPC, Fermilab from Oct. 2006 to June 2007
 - In order to validate new release of simulation and reconstruction package of CMSSW large statistics of Monte Carlo sample was generated.
- Conducted the simulation workshop for CMS
 - Participants included post-doctoral fellows, graduate students, system managers and software experts.
 - Learned the installation of CMS software and their use in physics analysis
- System administrator of Delhi group

List of Publications



- Abstract accepted in CHEP 2012
 - Enabling data analysis la PROOF on the Italian ATLAS-Tier2's using PoD
 - The ATLAS Computing activities and developments of the Italian Cloud
- Grid related publication
 - Multicore in Production: Advantages and Limits of the Multi-process Approach
ACAT, September 5-9, 2011, Uxbridge, London
 - Data analysis with GANGA: Accepted for publication in J.Phys.Conf.Series.
 - Distributed analysis functional testing using GangaRobot in the ATLAS experiment:
Accepted for publication in J.Phys.Conf.Series
 - Computing infrastructure for ATLAS data analysis in the Italian cloud: Accepted for
publication in J.Phys.Conf.Series
 - ATLAS Muon Calibration Frameowrk. Accepted for publication in
J.Phys.Conf.Series
 - A new CDF model for data movement based on SRM". M.K. Jha, ..., Doug
Benjamin, et al, Published in: J.Phys.Conf.Ser.219:062052,2010

Thanks !

Backup Slides

HammerCloud Web UI



Hammercloud | ATLAS

You are connected as gangarbt, [click here to Logout!](#)

Home	Clouds	Tests	HC Stats	Ganga Robot	Panda Dashb.	Administration
------	--------	-------	----------	-------------	--------------	----------------

Welcome to HammerCloud-ATLAS. Click "HC Stats" above to see the currently running jobs.

Running and Scheduled Stress Tests

state	id	host	template	start time (CET)	end time (CET)	clouds	sites	subm jobs	run jobs	comp jobs	fail jobs	tot jobs
running	10000340	voatlas73	Muon 15.6.6 PANDA default data-access	2010-07-07 09:03:00	2010-07-07 13:03:00	US	ANALY_MWT2, ANALY_MWT2_X, ANALY_SLAC, 2 more...	99	744	977	94	1914

Running and Scheduled Functional Tests

state	id	host	template	start time (CET)	end time (CET)	clouds	sites	subm jobs	run jobs	comp jobs	fail jobs	tot jobs
running	10000342	voatlas49	UA 15.6.9 Panda Filestager	2010-07-07 10:38:07	2010-07-08 10:38:07	CA_PANDA, DE_PANDA, ES_PANDA, 9 more...	ANALY_AGLT2, ANALY_ALBERTA, ANALY_ANLASC, 89 more...	58	60	258	11	405
running	10000338	voatlas73	D3PDMaker 15.6.10.6 PANDA default data-access Frontier/Squid test	2010-07-06 15:47:03	2010-07-07 15:47:03	CA_PANDA, DE_PANDA, ES_PANDA, 9 more...	ANALY_AGLT2, ANALY_ALBERTA, ANALY_ANLASC, 89 more...	156	32	567	167	941
running	10000337	voatlas73	D3PDMaker 15.6.10.6 LCG Frontier/Squid test	2010-07-06 15:41:03	2010-07-07 15:41:03	CA, DE, ES, 7 more...	ALBERTA-LCG2_MCDISK, AUSTRALIA-ATLAS_MCDISK, BEIJING-LCG2_MCDISK, 98 more...	8	26	445	183	662
running	10000336	voatlas49	UA 15.6.9 Panda	2010-07-06 13:12:04	2010-07-07 13:12:04	CA_PANDA, DE_PANDA, ES_PANDA, 9 more...	ANALY_AGLT2, ANALY_ALBERTA, ANALY_ANLASC, 89 more...	64	54	5366	150	5649
running	10000335	voatlas73	UA 15.6.9 LCG DQ2 Local	2010-07-06 13:12:04	2010-07-07 13:12:04	CA, DE, ES, 7 more...	ALBERTA-LCG2_MCDISK, AUSTRALIA-ATLAS_MCDISK, BEIJING-LCG2_MCDISK, 98 more...	19	44	2001	174	2240
running	10000334	voatlas73	UA 15.6.9 LCG DQ2 Filestager	2010-07-06 13:11:04	2010-07-07 13:11:04	CA, DE, ES, 7 more...	ALBERTA-LCG2_MCDISK, AUSTRALIA-ATLAS_MCDISK, BEIJING-LCG2_MCDISK, 98 more...	22	44	2073	116	2257

<http://hammercloud.cern.ch/atlas/>

Manoj K. Jha

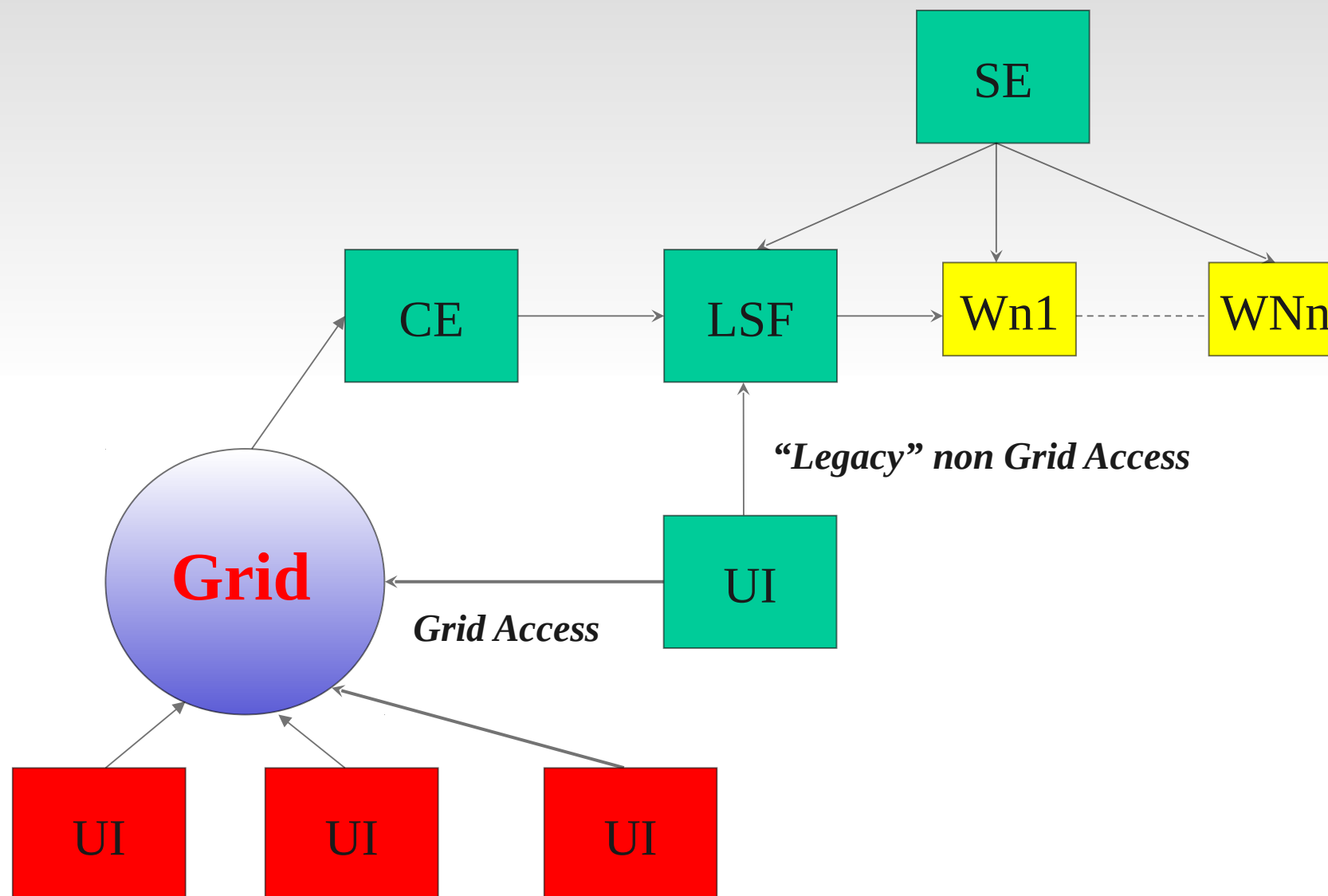
ATLAS Analysis in a Nutshell

- Data
 - Centrally organized data distribution by data management system (DQ2) according to computing model
- Experiment software (Athena) distribution kits
 - Centrally organized installation on EGEE, OSG and NG
 - Sites are moving toward CVMFS for availing software distribution kits on worker nodes
- User jobs
 - Model: “Job goes to data”
 - Tools for user job management: Ganga and Panda clients
- User output
 - Store on site scratchdisk or transfer on demand to remote disk
 - Retrieve output with DQ2 command line tools to local computer

Farming

- Tasks
 - Installation & management of Tier1 WNs and servers
 - Using Quattor (still some legacy lcfgng nodes around)
 - Deployment & configuration of OS & LCG middleware
 - HW maintenance management
 - Management of batch scheduler (LSF, torque)

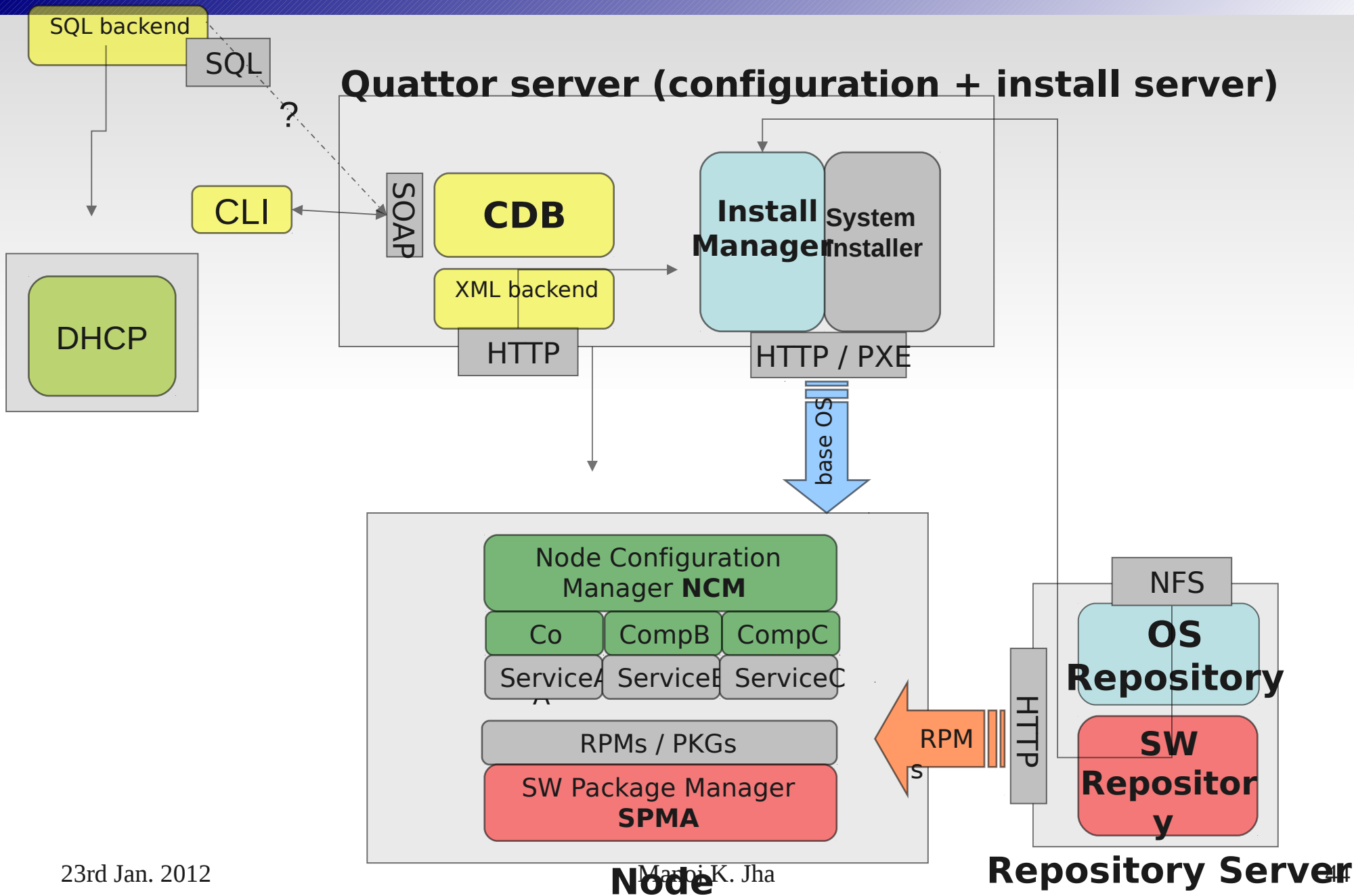
Access to Batch system



What we do with quattor ?

- Base OS installation
- Installation of different types of farm
 - LCG
 - Experiment specific farms
- We use quattor to keep updated the farm in terms of configuration and software.

Quattor architecture at CNAF



Components used



- grub
- nfs
- ldconf
- accounts
- authconfig
- afs
- ntp
- chkconfig
- altlogrotate
- cron
- globuscfg
- cmnconfig
- rm
- dirperm
- filecopy
- profile
- edglcg
- rgmaclient
- gridmapdir
- gsissh

Node installation process



1. Update local DB with node info
 - ✓ S/N, location, HW, Network, ecc...
2. DNS and DHCP automatically updated by DB update process
3. Update by hand `pro_site_databases.tpl`
 - ✓ `escape("wn-03-02-01-a.cr.cnaf.infn.it"),"131.154.192.151",`
 - ✓ `escape("wn-03-02-0a.cr.cnaf.infn.it"),"pro_hardware_machine_sun",`
4. Create and add to CDB the node profile
 - ✓ `cdb-simple-cli - -add profile_wn-03-02-01-a.tpl`
5. Configure PXE and KickStart for node
 - ✓ `aii-shellfe - -configure wn-03-02-01-a`
 - ✓ `aii-shellfe - -install wn-03-02-01-a`
6. Booting node (configure the correct boot device sequence)
7. DHCP supplies IP and location of kernel and KickStart configuration
8. AII takes care of installing and configuring the node
 - ✓ Installing Base OS
 - ✓ Reboot and execution of `ks-post-reboot` script
 - ✓ Install the Quattor client
 - ✓ Upgrade the system if required (`lsg, ...`)